# Assignment-based Subjective Questions

Q1 :- From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable?

Ans :- When analyzing the bike-sharing dataset to understand the factors that impact bike demand:

Season: Bike rentals are higher in summer and fall, likely because the weather is nicer, making biking more enjoyable.

Yearly Trend: There was a noticeable increase in bike usage from 2018 to 2019, indicating that more people are using the bike-sharing service over time.

Holiday Impact: Less demand for bikes on holidays, probably because fewer people are commuting to work or running errands.

Working Days: More bikes are rented on working days, suggesting that people often use them to commute to work.

Weather Conditions: Clear or partly cloudy days see the highest bike demand, while rainy or snowy days see the least, as bad weather makes biking less appealing and safe.

These observations were derived from the patterns seen in the data, particularly from box plots that show how bike demand varies with different conditions.

Q2 :- Why is it important to use drop_first=True during dummy variable creation?

Ans :- Using drop_first=True during dummy variable creation helps prevent a common problem called "multicollinearity" when working with statistical models. Multicollinearity occurs when one variable can be perfectly predicted from the others, making it hard to determine the individual effects of each variable on the outcome.

In simple terms, when you create dummy variables (which convert categorical data into a numerical format that models can use), if you include a column for every category, you are adding redundant information. For instance, if you have a variable for "Season" (Spring, Summer, Autumn, Winter) and create dummies for all four, knowing the values for three of the seasons (e.g., if Spring, Summer, and Autumn are all 0), you automatically know the value of the fourth season (Winter must be 1). This redundancy can skew or confuse the analysis.

By using drop_first=True, you drop the first dummy variable (e.g., Spring) and keep the rest (Summer, Autumn, Winter). You're not losing information but simplifying the model, as the dropped category can be inferred from the others (if Summer, Autumn, and Winter are all 0, it implies that it is Spring). This approach helps ensure that the model is more stable and provides more reliable insights.

Q3 :- Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable?

Ans :- Temperature has the highest correlation with target Variable with around 0.63

Q4 :- How did you validate the assumptions of Linear Regression after building the model on the training set?

Ans :- After building a linear regression model on the training set, validating the assumptions of the model ensures its reliability and accuracy. Here are key assumptions of linear regression and methods to validate them:

Linearity: This assumption states that there is a linear relationship between the predictors and the dependent variable.

Validation: To check this, you can plot the observed vs. predicted values or use residuals vs. fitted values plots. If the relationship is linear, points should roughly form a straight line. Non-linear patterns suggest transforming the variables or using polynomial terms.

Normality of Residuals: Residuals (the differences between observed and predicted values) should be normally distributed.

Validation: A histogram or a Q-Q plot (quantile-quantile plot) of the residuals can be used to assess normality. If the residuals follow a straight line on the Q-Q plot, they are likely normally distributed.

Homoscedasticity: This implies that the residuals have constant variance at all levels of the independent variable.

Validation: Plot residuals against fitted values. The plot should not show any patterns or funnel shape, indicating increasing or decreasing error variance.

Independence of Residuals: There should be no correlation among residuals.

Validation: Examine a residual autocorrelation plot (e.g., Durbin-Watson test can be used). The residuals should not display correlation.

No Multicollinearity: Predictors should not be perfectly multicollinear. In other words, one variable should not be too closely linearly related to another.

Validation: Check Variance Inflation Factor (VIF); a VIF above 5 (or 10, based on the stringent level you choose) suggests multicollinearity and that variable might need to be removed or adjusted.

No significant outliers, high-leverage points, or influential points: Outliers and influential points can disproportionately impact the fitting of the model.

Validation: Leverage plots or influence plots can help detect these points. Removing or adjusting outliers might be necessary to improve the model.

Each of these validations helps in enhancing the model by confirming each assumption or indicating necessary adjustments, such as transformations, removing variables, or treating outliers. Such scrutiny helps in building a robust linear regression model that provides reliable predictions.

Q5 :- Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes?

Fall season :-  it associated with a higher demand for bikes compared to otherseasons, probably because of favorable weather conditions.

Temperature  :- has a strong positive effect on bike demand; as temperatures rise,more bikes are rented.

Year 2019 :-  shows a significant increase in bike rentals compared to 2018,reflecting overall growth in the bike-sharing system.

# General Subjective Questions

Q1 :- Explain the linear regression algorithm in detail.

Linear regression is a statistical approach employing a linear equation to exemplify the relationship between a dependent variable (also known as the target) and one or more independent variables (referred to as predictors). It is straightforward and commonly used for predictive modeling purposes.

Objectives:

Prediction: Estimating the values of the dependent variable using the independent variables.

Exploration of Relationships: Assessing the magnitude and type of association between the dependent and independent variables.

Key Components:

Dependent Variable (Y): This is the variable whose values are predicted.

Independent Variables (X1, X2, ..., Xn): These are the predictors or input variables used in the prediction process.

Linear Equation Representation: The model asserts the relationship through the equation:

$$ Y = \beta_0 + \beta_1X_1 + \beta_2X_2 + ... + \beta_nX_n + \epsilon $$

$Y$: Dependent variable

$X_1, X_2, ..., X_n$: Independent variables

$\beta_0$: Y-intercept

$\beta_1, \beta_2, ..., \beta_n$: Coefficients representing weights of the independent variables

$\epsilon$: Error term

Steps in Linear Regression:

Data Collection: Compile data that includes both dependent and independent variables.

Exploratory Data Analysis (EDA): Use statistical summaries and visualizations to gain insights into the data.

Data Preprocessing: Address any missing values, encode categorical variables, and standardize or normalize data as required.

Model Construction: Apply the linear regression model to the collected training data.

Model Assessment: Evaluate the model's effectiveness using statistical measures like R-squared, Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE).

Residual Examination: Analyze the residuals to ensure the model adheres to the essential assumptions of linear regression.

Prediction: Utilize the model for forecasting values based on new or previously unseen data.

Assumptions in Linear Regression:

Linearity: There exists a linear correlation between all independent and the dependent variables.

Independence: The residuals from the predictions are independent from each other.

Homoscedasticity: The residuals possess constant variance at varying levels of the independent variable(s).

Normality of Residuals: The distribution of residuals should approximate normality.

Adhering to these assumptions and thoroughly validating the model, linear regression can serve as a powerful and interpretable tool for elucidating the relationship between a dependent variable and multiple independent variables. It can deliver meaningful insights and precise predictions essential for informed decision-making.

Q2 :- Explain the Anscombe's quartet in detail.

Ans :- Anscombe's quartet is a set of four different datasets devised by the statistician Francis Anscombe in 1973 to demonstrate the importance of visualizing data before analysing it with statistical models, particularly linear regression. Each of the four datasets comprises eleven points (x, y) and, at first glance, they have nearly identical statistical properties, misleading an analyst who relies solely on summary statistics.

Key Identical Statistics in Each Dataset:

Mean of x: The mean of the x values is approximately 9 for all datasets.

Mean of y: The mean of the y values is approximately 7.50 for all datasets.

Variance of x: The variance of x is approximately 11.

Variance of y: The variance of y is approximately 4.12.

Correlation between x and y: The correlation coefficient is approximately 0.816 for all datasets.

Linear Regression Line: When a linear regression is performed on each dataset, the line of best fit is approximately y = 3.00 + 0.500x.

Coefficient of Determination (R-squared): All datasets exhibit a similar R-squared value, showing a high level of explanation of variability.

Despite these similarities in key statistical summaries, the distribution and relationship patterns of the data points in each quartet are distinctively different when graphed. This highlights different underlying relationships and structures:

Details of Each Dataset:

Dataset I: Features a simple linear relationship between x and y. It corresponds closely to the assumptions of linear regression.

Dataset II: Shows a curvature (quadratic relationship). This dataset veers off linear behavior suggesting that a simple linear model might not be appropriate.

Dataset III: Contains an outlier which influences the slope of the regression line significantly. Without this outlier, the correlation between x and y would be substantially weaker.

Dataset IV: Demonstrates how a single influential point can also drastically affect the regression line. It consists of x-values that are mostly constants except for one outlier, which heavily influences the regression equation.

Importance Highlighted by Anscombe's Quartet:

Visual Inspection: The quartet underscores the necessity of plotting your data before relying on statistical metrics. Visualizing data can reveal patterns, anomalies, or deviations that pure statistical analysis might fail to detect.

Limitation of Statistics Alone: It shows that similar statistical properties can lead to incorrect assumptions about the underlying data structure or relationship. This serves as a caution against using only summary statistics to describe or analyze data.

Robustness of Statistical Analysis: Anscombe's Quartet promotes the idea of applying robust statistical analysis that involves checking the data's behavior through visualization in addition to performing standard statistical tests.

The Anscombe's quartet is often used in statistics, data analysis, and data science education to emphasize the importance of data visualization alongside numerical analysis. It serves as a classic reminder that robust statistical analysis is both an art and a science, requiring a comprehensive approach beyond surface-level metrics.

Q3 :- What is Pearson's R?

Ans :-

Pearson's R, also known as the Pearson product-moment correlation coefficient (PPMCC) or simply the Pearson correlation coefficient, is a statistical measure that evaluates the linear relationship between two quantitative variables. Named after its developer, Karl Pearson, this coefficient provides a value between -1 and 1 that indicates the strength and direction of the linear association between the variables.

Calculation of Pearson's R:

Pearson's R is calculated using the formula:

$$ r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}} $$

where:

( $n$ ) is the number of data points,

( $\sum xy$ ) is the sum of the product of paired scores,

( $\sum x$ ) and ( $\sum y$ ) are the sums of the ( $x$ ) scores and ( $y$ ) scores respectively,

( $\sum x^2$ ) and ( $\sum y^2$ ) are the sums of the squares of the ( $x$ ) and ( $y$ ) scores.

Interpretation of Pearson's R:

1 or -1: A correlation of ±1 indicates a perfect positive (1) or negative (-1) linear relationship between the variables. All data points lie exactly on a line with positive or negative slope.

0: A correlation of 0 suggests no linear relationship between the variables. However, it does not imply the absence of any relationship; variables may have other non-linear associations.

Between 0 and ±1: Values that are not exactly -1, 0, or 1 denote degrees of linear relationships. Closer to ±1, the stronger the correlation, with signs indicating the direction (positive or negative).

Significance:

Positive values indicate a positive correlation, where increases in one variable correspond with increases in the second variable. Conversely, negative values suggest a negative correlation, where increases in one variable correspond with decreases in the other.

Usage:

Pearson's correlation coefficient is widely used in the sciences and economics to calculate correlations between variables. It helps in data analysis for:

Preliminary investigations of relationships.

Providing insights in conjuncture with other analyses such as regression, causal modeling, and prediction.

Validation of relationships and dependencies between quantities.

Limitations:

While Pearson's R is extremely useful, it has some limitations:

It only measures linear relationships and fails to detect non-linear relationships.

It can be affected significantly by outliers.

A statistically significant Pearson correlation does not imply causality between the variables.

Conclusion:

Pearson's R is a foundational tool in statistical analysis for exploring and quantifying linear relationships between pairs of continuous variables. Its simplicity and interpretability make it an essential part of the data analyst's toolkit, though care must be taken to consider its assumptions and limitations when interpreting results.

Q4 :- What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling?

Ans :- What is Scaling?

Scaling is a method used in data preprocessing to standardize the range of independent variables or features of data. In other words, it involves transforming the values of numeric columns in the dataset to a common scale, without distorting differences in the ranges of values or losing information.

Why is Scaling Performed?

Scaling is crucial because many machine learning algorithms, like logistic regression, linear regression, and neural networks, perform better when input numerical variables fall within a similar scale. This is particularly true for algorithms that use distance calculations, as having features on the same scale allows the model to converge faster and appropriately weigh all the features. For example, one feature

measured in kilometers and another in grammes could disproportionately influence the model due to their differing scales.

Normalized Scaling vs. Standardized Scaling:

Normalized Scaling (Min-Max Scaling): This type of scaling rescales the feature to a fixed range, usually 0 to 1, or -1 to 1, using the formula:

$$X_{\text{norm}} = \frac{X - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}}$$

Min-Max scaling is useful when you need to preserve exactly zero points in the data, but it is sensitive to outliers.

Standardized Scaling (Z-score Normalization): This method involves rescaling the data to have a mean (μ) of 0 and a standard deviation (σ) of 1. The formula used is:

$$X_{\text{std}} = \frac{X - \mu}{\sigma}$$

Standardized scaling is less affected by outliers and often used if the data does not have a specific range to anchor to.

Q5 :- You might have observed that sometimes the value of VIF is infinite. Why does this happen?

Ans :- Infinite Value of VIF:

The Variance Inflation Factor (VIF) quantifies how much the variance of an estimated regression coefficient increases if your predictors are correlated. If predictors are uncorrelated, VIF will be 1. VIF provides an index that measures how much the variance of an estimated regression coefficient is increased because of collinearity.

An infinite VIF value usually arises when there is perfect multicollinearity between independent variables, meaning one variable can be perfectly linearly predicted by the others with absolute accuracy. This situation destroys the ability to estimate the coefficients of the predictors, as it becomes challenging to determine where one effect starts and another ends.

Q6 :- What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression.

Ans :- A Q-Q (quantile-quantile) plot is a graphical tool to help us assess if a set of data plausibly came from some theoretical distribution such as a Normal or exponential. It compares two probability distributions by plotting their quantiles against each other.

Use and Importance in Linear Regression:

In linear regression, assumptions regarding the normality of residuals are crucial as they influence the validity of the hypothesis tests and confidence intervals calculated for the coefficients. The Q-Q plot is particularly used to assess the normality of residuals. A linear pattern in the Q-Q plot indicates that the residuals are normally distributed. Deviations from this line suggest departures from normality, signaling potential issues with inferences from standard linear regression models.

Using a Q-Q plot allows researchers and analysts to visually check the assumption of normality, which is a cornerstone for many parametric methods in statistics and critical for the accuracy of many linear regression models.